
Lecture 2



Logistics and announcement

- CVN has confirmed that class will be recorded
- After hours office hours will be available by appointment only
 - Please email Hongyi in case you need to schedule it
- Due to popular demand, final exam will be in class on May 4
- Programming homework 1 is coming out
 - Details on website
 - You should all receive coupons for Google Cloud
 - Please use your Columbia emails to log in
- We will reschedule Feb 10 class

Recap of lecture 1

■ Latency vs. throughput

■ Latency rules of thumb:

- Memory access: 100ns
- Read a small object within the same region in a data center: 100,000ns, 100us
- Run a SQL query on a flash database: 1,000,000ns, 1ms
- Run a SQL query on a disk database: 20,000,000ns, 20ms
- Roundtrip time over the internet: 100,000,000ns, 100ms

■ $Speedup(E) = \frac{ExTime_{old}}{ExTime_{new}} = \frac{1}{(1-p) + \frac{p}{f}}$

■ Amdahl's law: speedup bounded by:

■ Amdahl's law: speedup bounded by:

- $\frac{1}{fraction\ of\ time\ not\ enhanced}$
- In other words, make the common case fast!

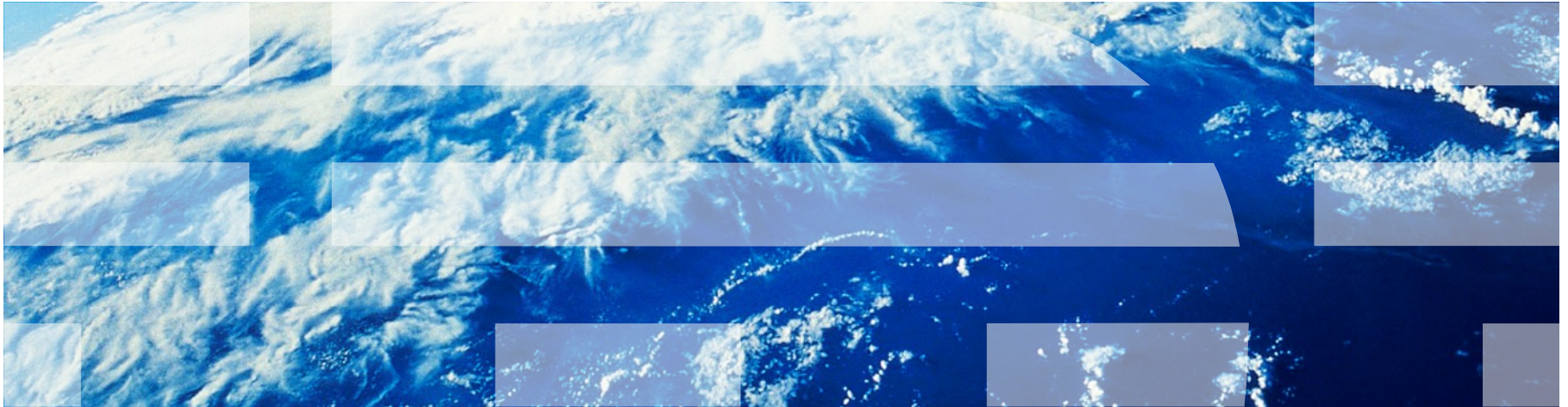
- In other words, make the common case fast!

Today's class

- Data centers
 - Example: Google
 - Racks/Slots/Rows
 - Networking
 - Power
 - Failures
- Relational model

Adapted from Mendel Rosenblum and Jeff Dean

The Infrastructure of Big Data



Motivating example: Google web search (1999 vs. 2010)

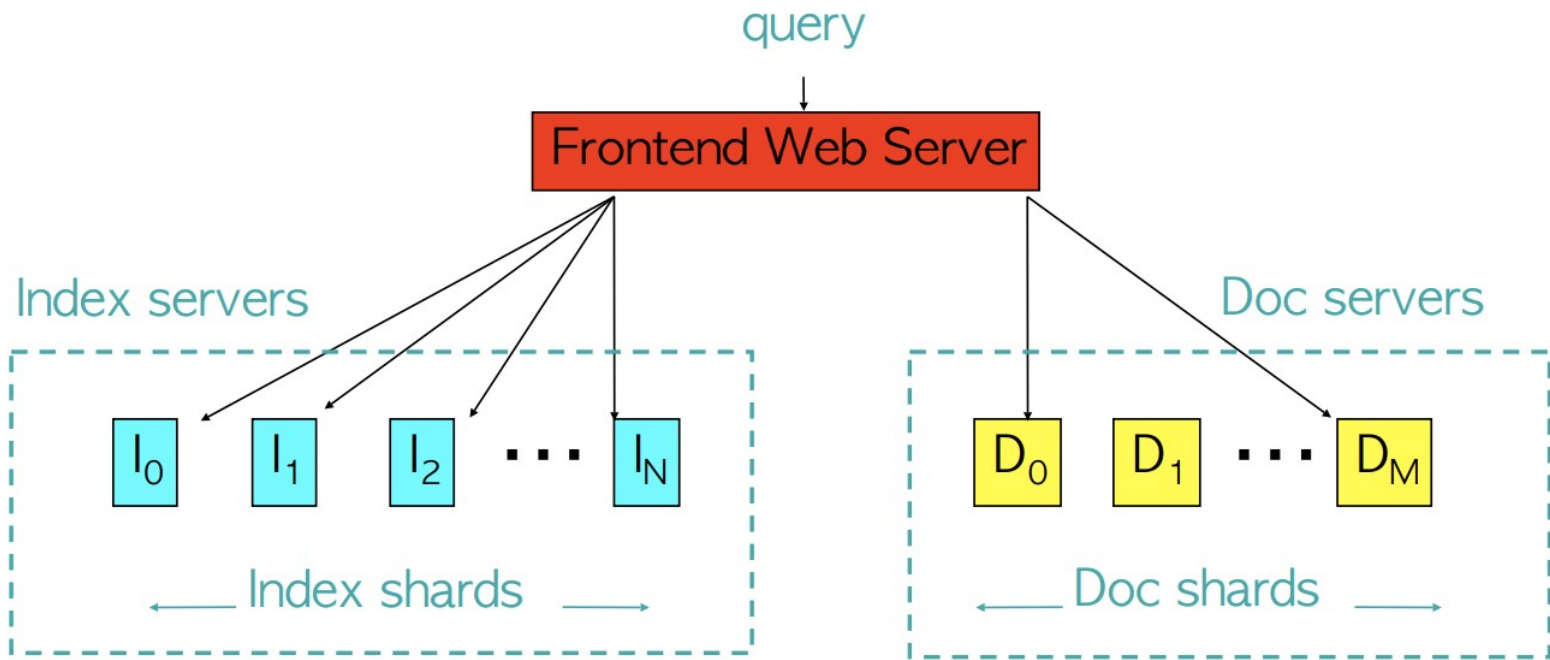
- # docs: tens of millions to tens of billions ~1000X
- Queries processed/day: ~1000X
- Per doc info in index: ~3X
- Update latency: months to tens of seconds ~50000X
- Average query latency: 1 seconds to 0.2 seconds ~5X

- More machines * faster machines: ~1000X

Google Circa 1997 (definitely not big data)



Google infrastructure circa 1997 could fit in a single room



Scaling up

- What happens when a server doesn't fit in a single room?
- What happens if we need 1000X more servers?
- The cloud to the rescue!
 - Also known as... **data centers**

Evolution of data centers

- 🐦 1960's, 1970's: a few very large time-shared computers
- 🐦 1980's, 1990's: heterogeneous collection of lots of smaller machines.
- 🐦 Today and into the future:
 - Data centers contain large numbers of nearly identical machines
 - Geographically spread around the world
 - Individual applications can use thousands of machines simultaneously
- 🐦 Companies consider data center technology a trade-secret
 - Limited public discussion of the state of the art from industry leaders

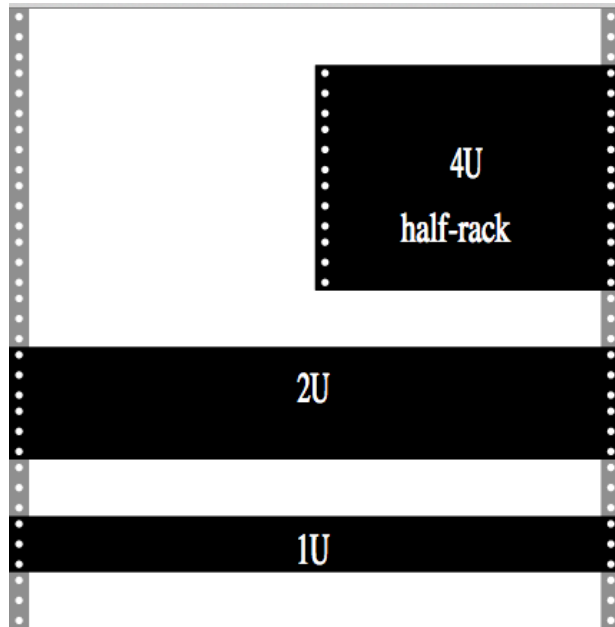
Typical specs for a data center today

- 🐦 15-40 megawatts power (Limiting factor)
- 🐦 50,000-200,000 servers
- 🐦 \$1B construction cost
- 🐦 Onsite staff (security, administration): 15

Rack

- ✦ Typically is 19 or 23 inches wide
- ✦ Typically 42 U
 - U or RU is a Rack Unit - 1.75 inches

✦ Slots:



Rock Slots

✚ Slots hold power distribution, servers, storage, networking equipment

✚ Typical server: 2U

- 8-128 cores
- DRAM: 32-512 GB

✚ Typical storage: 2U

- ✚ 30 drives

✚ Typical Network: 1U

- 72 10GB



Row/Cluster

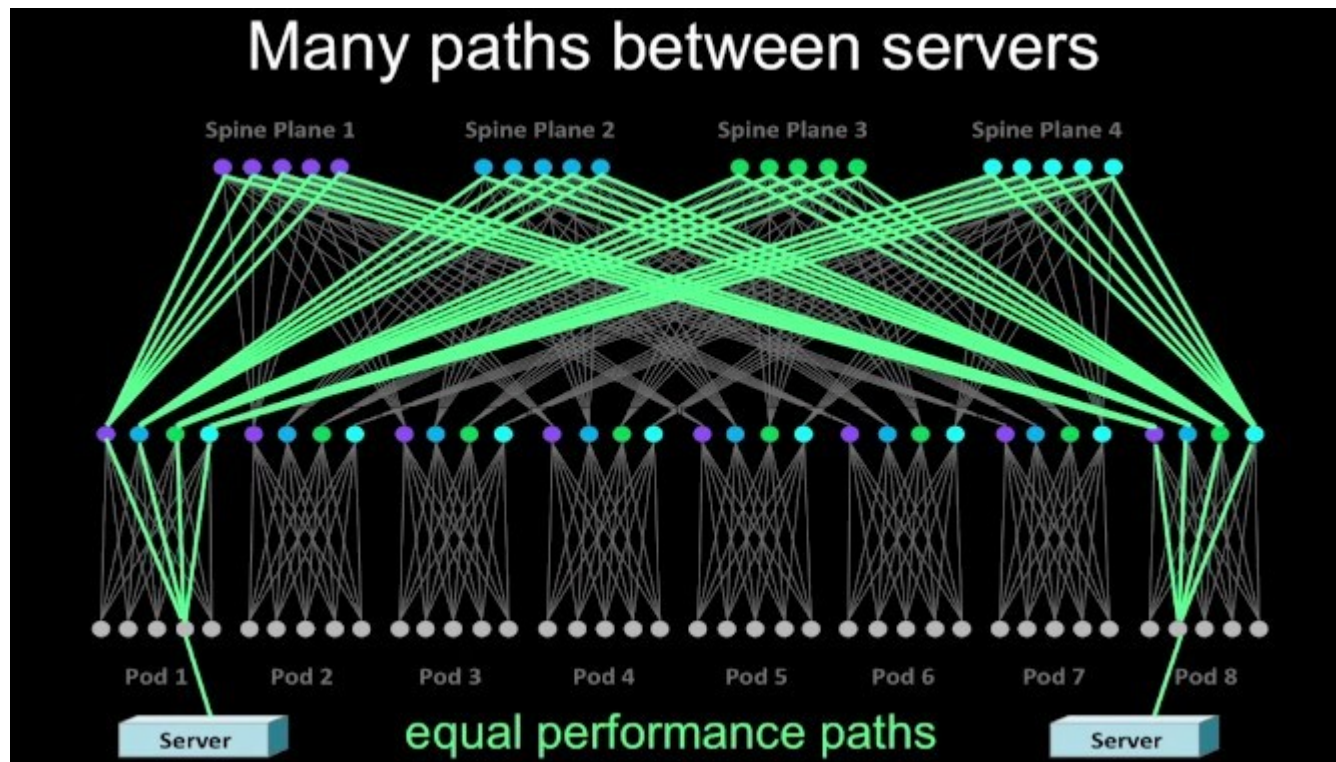
👤 30+ racks



Networking - Switch locations

- 👉 Top-of-rack switch
 - Connecting machines in rack
 - Multiple links going to end-of-row routers
 - 👉 End-of-row router
 - Aggregate row of machines
 - Multiple links going to core routers
 - 👉 Core router
 - Multiple core routers
- 👉 Each of these have different latencies, throughput

Multipath routing



Ideal: "full bisection bandwidth"

🐼 Would like network where everyone has a private channel to everyone else

- (cross-bar topology)
- Why is this useful?

🐼 In practice today:

- Assumes applications have locality to rack or row but this is hard to achieve in practice.
- Some datacenter networking problems are fundamental: Two machines transferring to the same machines
 - When could this happen?

Power Usage Effectiveness (PUE)

- 🐦 Early data centers built with off-the-shelf components
 - Standard servers
 - HVAC unit designs from malls

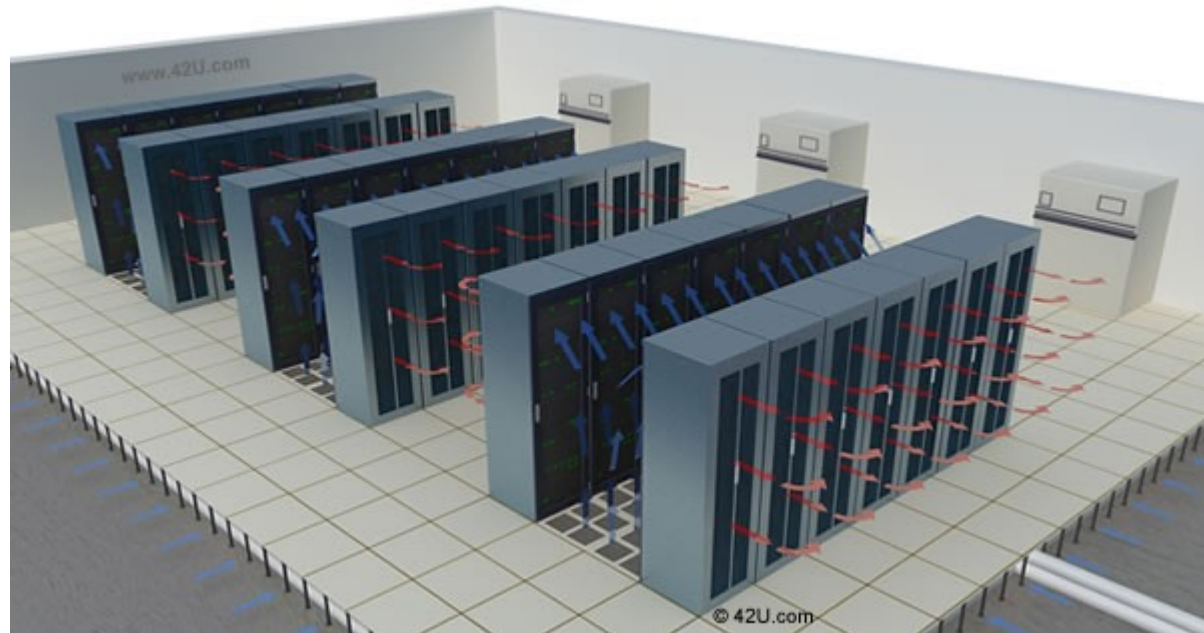
- 🐦 Inefficient: Early data centers had PUE of 1.7-2.0

$$\text{PUE ratio} = \frac{\text{Total Facility Power}}{\text{Server/Network Power}}$$

- 🐦 Best-published number (Facebook): 1.07 (no air-conditioning!)
- 🐦 Power is about 25% of monthly operating cost
 - And is a limiting factor in how large the datacenter can be

Energy Efficient Data Centers

- ✎ Better power distribution - Fewer transformers
- ✎ Better cooling - use environment (air/water) rather than air conditioning
 - Bring in outside air
 - Evaporate some water
- ✎ IT Equipment range
 - OK up to +115°F



Backup Power

- 🐦 Massive amount of batteries to tolerate short glitches in power
 - Just need long enough for backup generators to startup
- 🐦 How do glitches occur?
 - Thunder, earthquake, power loss from power company, cyber attack, ...
- 🐦 Massive collections of backup generators
- 🐦 Huge fuel tanks to provide fuel for the generators
- 🐦 Fuel replenishment transportation network (e.g. fuel trucks)

Fault Tolerance

- 👉 At the scale of new data centers, things are breaking constantly
- 👉 Every aspect of the data center must be able to tolerate failures
- 👉 Solution: Redundancy
 - Multiple independent copies of all data
 - Multiple independent network connections
 - Multiple copies of every services

Failures in first year for a new data center (Jeff Dean)

- ~thousands of **hard drive failures**
- ~1000 **individual machine failures**
- ~dozens of minor **30-second blips** for DNS
- ~3 **router failures** (have to immediately pull traffic for an hour)
- ~12 **router reloads** (takes out DNS and external VIPs for a couple minutes)
- ~8 **network maintenances** (4 might cause ~30-minute random connectivity losses)
- ~5 **racks go wonky** (40-80 machines see 50% packet loss)
- ~20 **rack failures** (40-80 machines instantly disappear, 1-6 hours to get back)
- ~1 **network rewiring** (rolling ~5% of machines down over 2-day span)
- ~1 **rack-move** (plenty of warning, ~500-1000 machines powered down, ~6 hours)
- ~1 **PDU failure** (~500-1000 machines suddenly disappear, ~6 hours to come back)
- ~0.5 **overheating** (power down most machines in <5 mins, ~1-2 days to recover)

□ **Reliability must come from software!**

Choose data center location drivers

- 👤 Plentiful, inexpensive electricity
 - Examples - Oregon: Hydroelectric; Iowa: Wind
- 👤 Good network connections
 - Access to the Internet backbone
- 👤 Inexpensive land
- 👤 Geographically near users
 - Speed of light latency
 - Country laws (e.g. Our citizen's data must be kept in our country.)
- 👤 Available labor pool

Google Data Centers

Americas

Berkeley County, South Carolina
Council Bluffs, Iowa
Douglas County, Georgia
Quilicura, Chile
Jackson County, Alabama
Mayes County, Oklahoma
Lenoir, North Carolina
The Dalles, Oregon

Asia

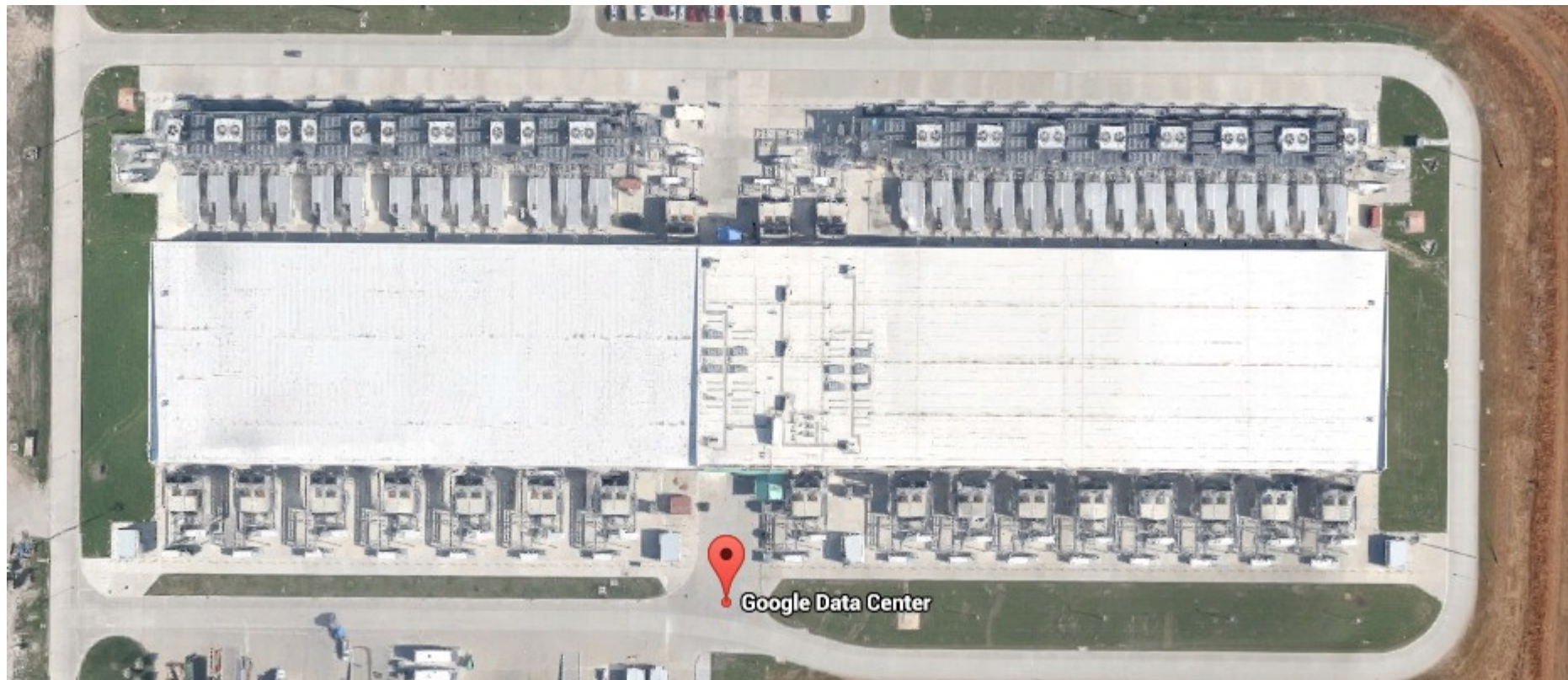
Changhua County, Taiwan
Singapore

Europe

Hamina, Finland
St Ghislain, Belgium
Dublin, Ireland
Eemshaven, Netherlands



Google Data Center - Council Bluffs, Iowa, USA



Google data center pictures: Council Bluffs

